# Machine Learning 1

Lecture 1.2 - What is Machine Learning?

*Erik Bekkers*

*(Bishop 1.0 and 1.1)*
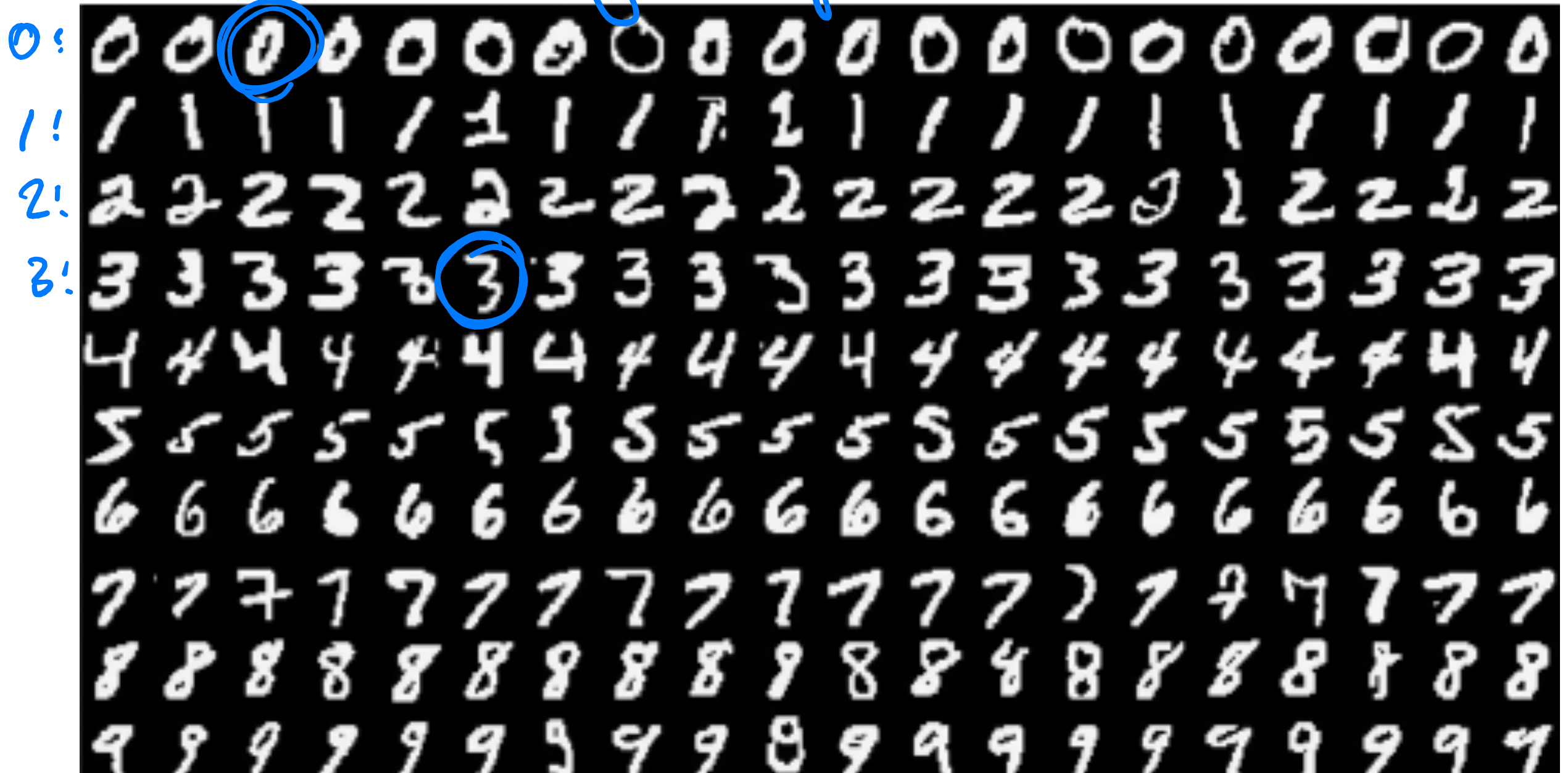
*Slide credits: Rianne van den Berg*

# What is machine learning?

"A computer program is said to learn from experience E with respect to some class of tasks T and performance measure  P if its performance at tasks in T, as measured by P, improves with experience E."
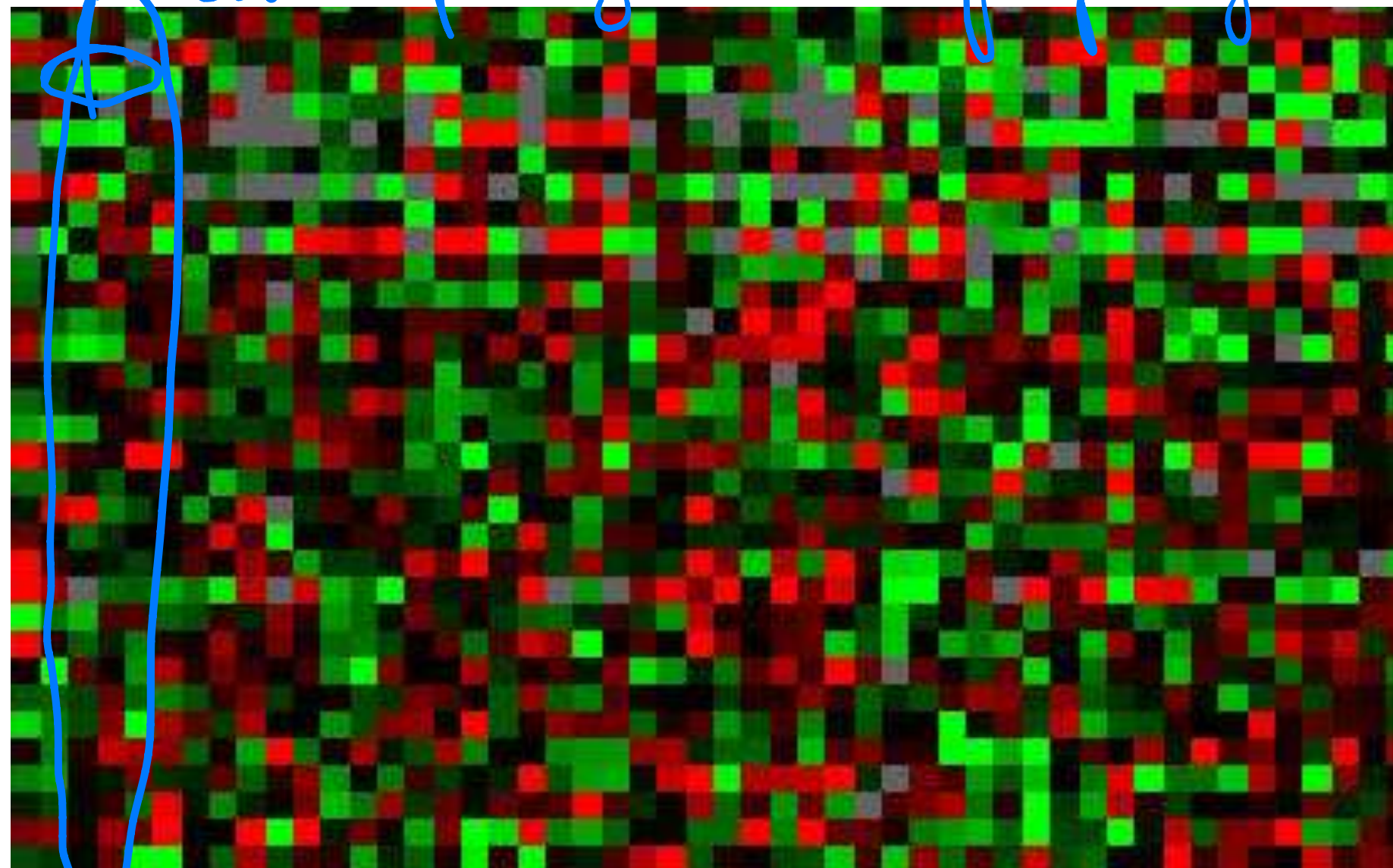
- Tom M. Mitchell

*Machine Learning, Tom Mitchell, McGraw Hill, 1997*

# E: Experience

Handwritten digits of size 28 × 28



MNIST dataset

# E: Experience



Expression matrix of genes (rows) for 64 human tumor samples (columns).  [source: ESL 1.3]

# E: Experience



Examples of spam emails. [source: Yesware]

# T: Class of tasks

Classification:



2 = 2

8 = 8

0 = 8

9 = 8 or 9

8 = 8 or 9

Spam
Spam
Spam
OK
Spam

**Time is running out. Save 50% on all the best moments!** - 50% Off Photo Purchase, $3.99 T

**you're so close to FREE snacks!** - we love you to try our delicious snacks | graze claim your

**Last Chance: Start 2016 with 50% off ___.com** - Get more from your 2016 with ___.com -

**Additional Incentive** - Great news Elise, from now through the end of the month ___ is offering

**PROOF: Diabetes Reversed 100% Naturally** - To receive this email in your inbox and activate t

All caps

Trigger phrase

Price
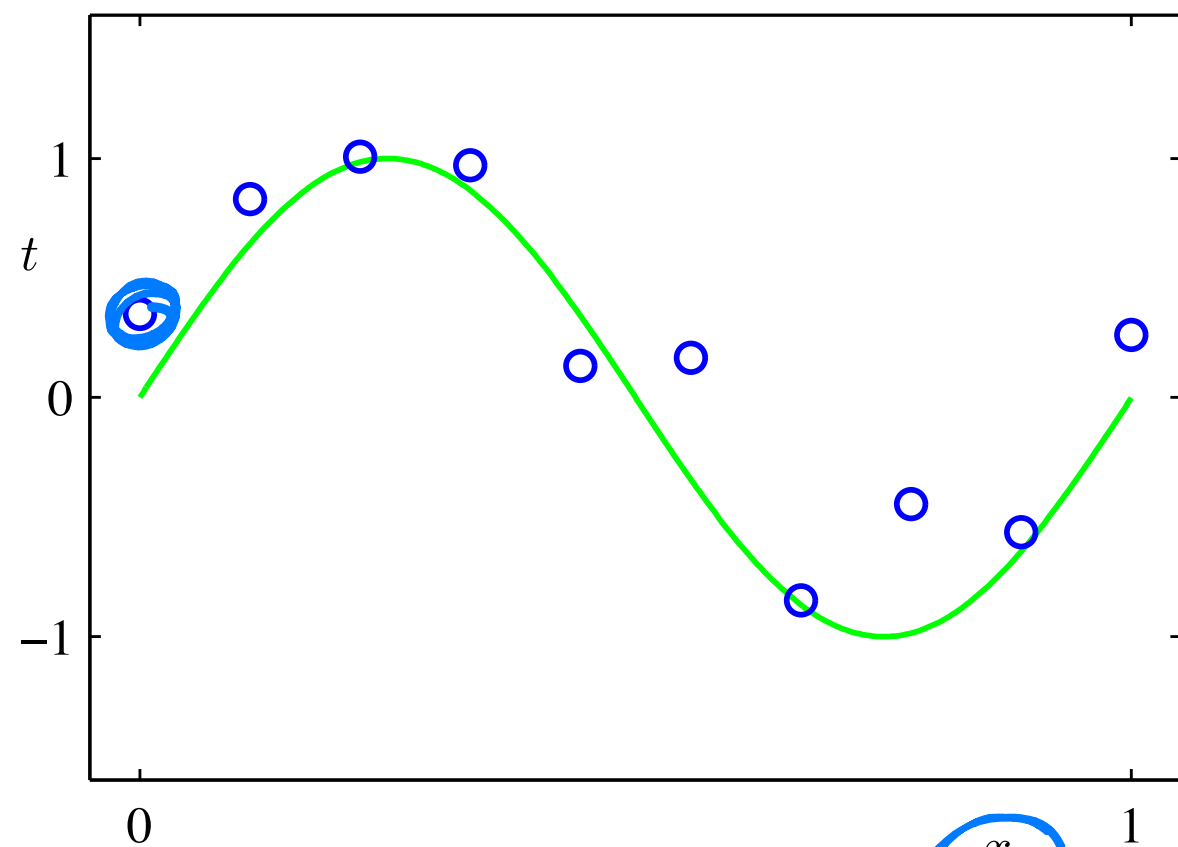
Exclamation point

Attachment

# T: Class of tasks

Regression

input :  $x$

target :  $t = \sin(2\pi x) + \varepsilon$

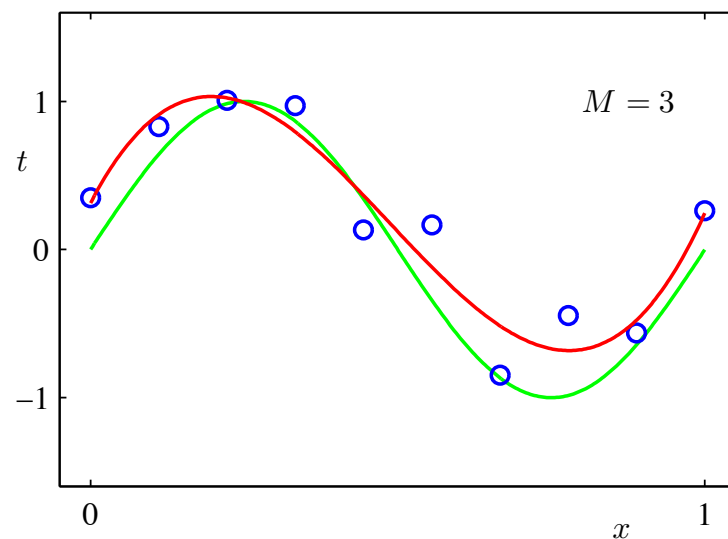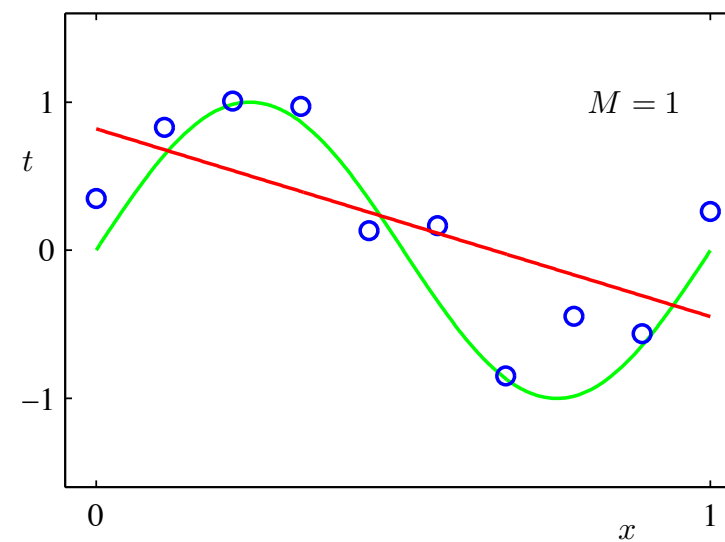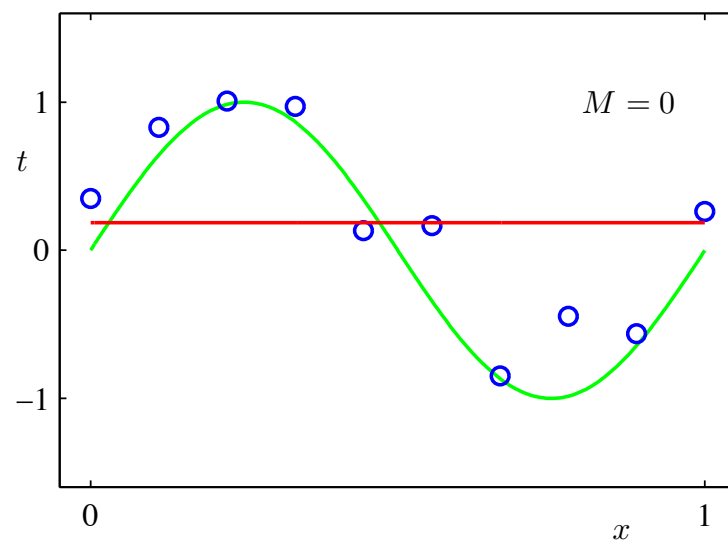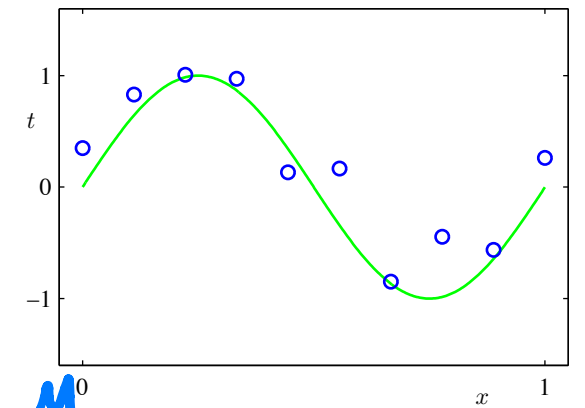noise :  $\varepsilon \sim \mathcal{N}(0, 1)$

# T: Class of tasks

Regression

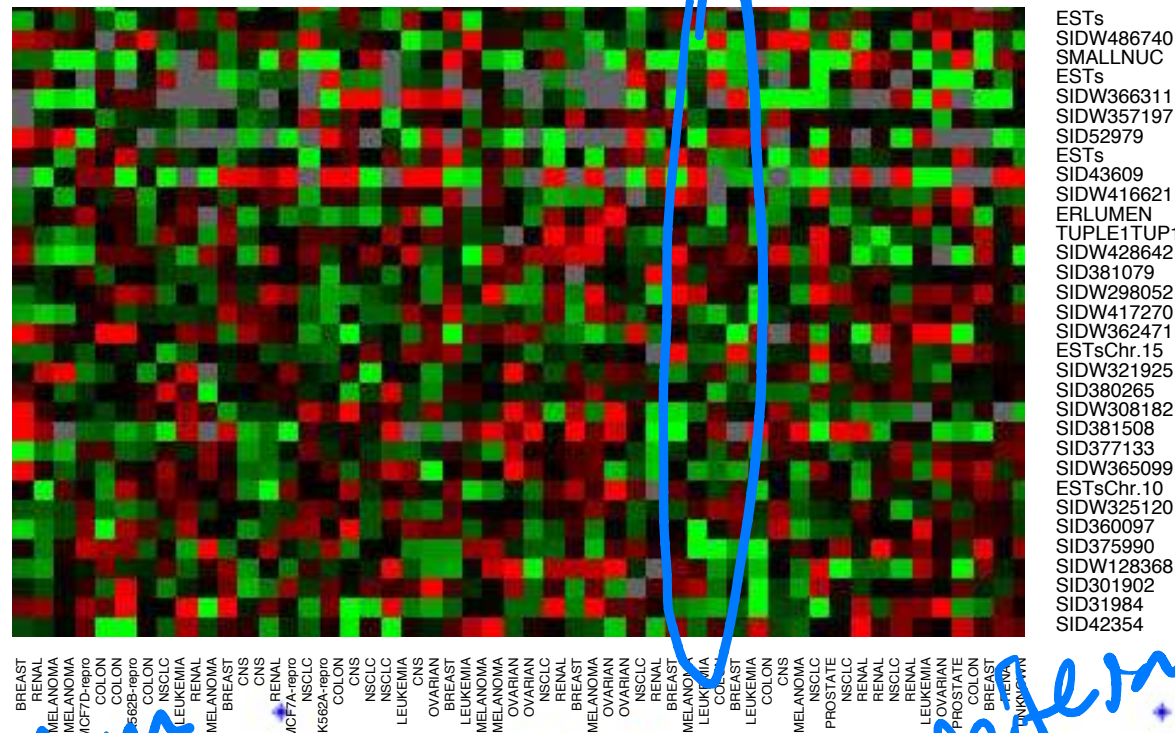$$f(x) = w_0 + w_1 x + w_2 x^2 + \cdots + w_M x^M$$



Polynomials of order $M$ (red) fit to data constructed as $t = \sin(2\pi x) + \varepsilon$ (green)

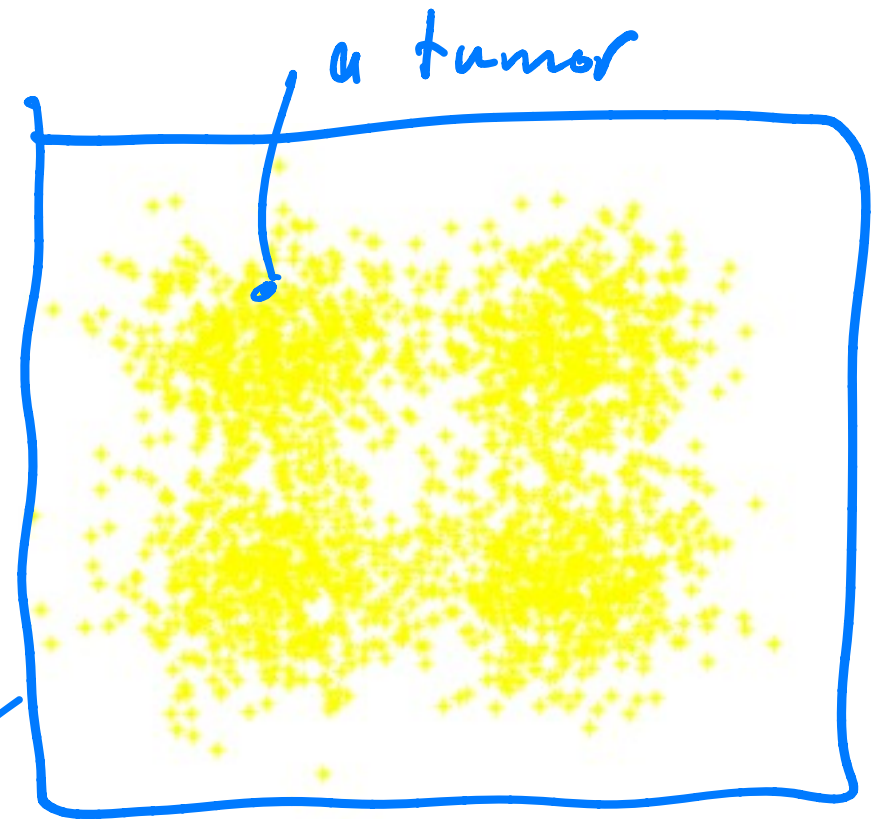# T: Class of tasks

Clustering

*N-dim Vector*

*2D point*

*a tumor*

*(most similar)*

*reference book for closest cluster*

*random assign each $x_i$ to one of 4 classes*

*$\mu_i$*

ESTs
SIDW486740
SMALLNUC
ESTs
SIDW366311
SIDW357197
SID52979
ESTs
SID43609
SIDW416621
ERLUMEN
TUPLE1TUP
SIDW428642
SID381079
SIDW298052
SIDW417270
SIDW362471
ESTsChr.15
SIDW321925
SID380265
SIDW308182
SID381508
SID377133
SIDW365099
ESTsChr.10
SIDW325120
SID360097
SID375990
SIDW128368
SID301902
SID31984
SID42354

Expression matrix of genes (rows) for 64 human tumor samples (columns). [source: ESL 1.3]
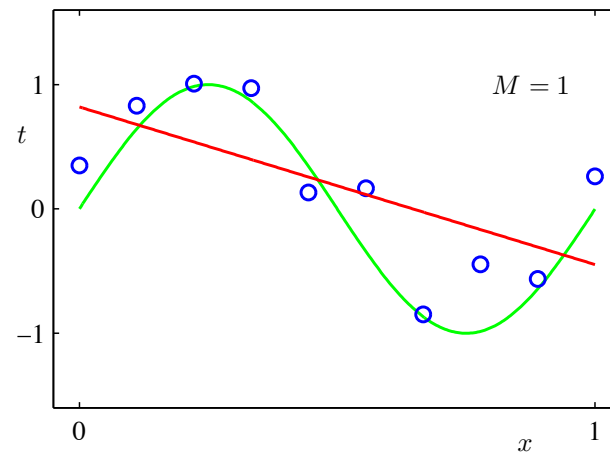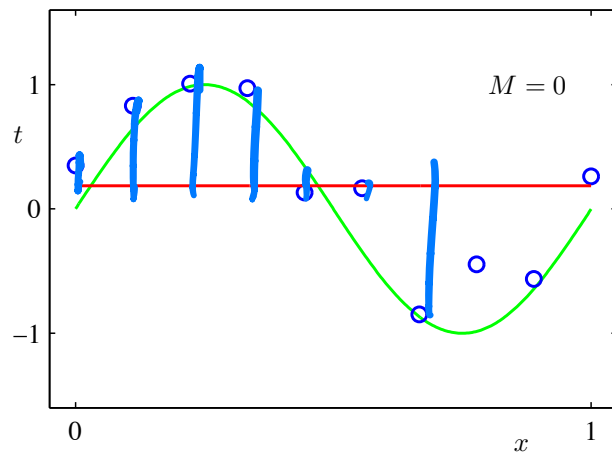
# P: Performance measure

Classification



$$\text{accuracy}(y, \hat{y}) = \frac{1}{n_{\text{samples}}} \sum_{i=1}^{n_{\text{samples}}} \mathbb{1}\left[y_i = \hat{y}_i\right]$$

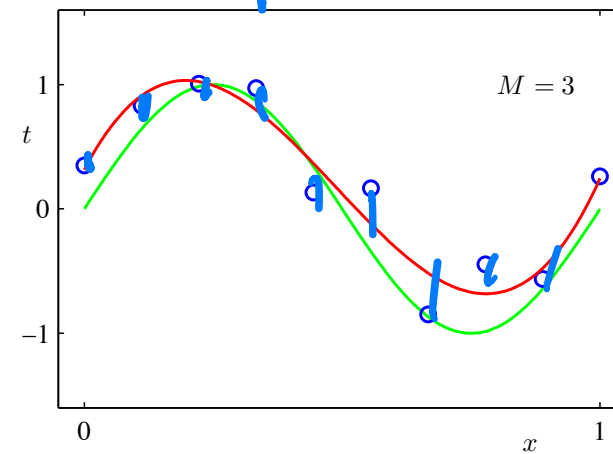$$\text{indicator funct.} = \begin{cases} 1 & \text{if } y_i = \hat{y}_i \\ 0 & \text{otherwise} \end{cases}$$
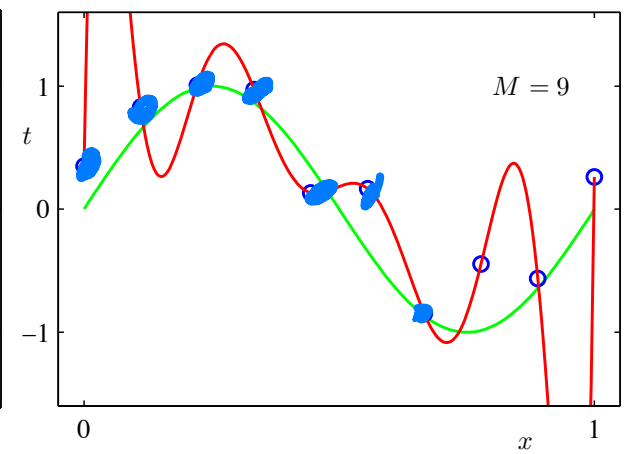
# P: Performance measure
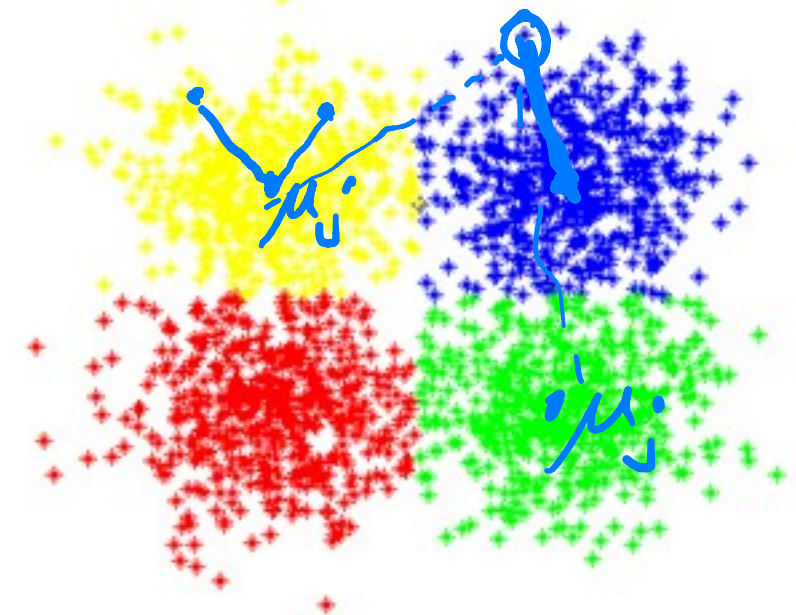
Regression

*large MSE*

*small MSE*

*MSE = 0*



Plots with $M = 0$, $M = 1$, $M = 3$, $M = 9$

*mean squared error*

$$\mathrm{MSE}(y, \hat{y}) = \frac{1}{n_{\mathrm{samples}}} \sum_{i=1}^{n_{\mathrm{samples}}} \left( y_i - \hat{y}_i \right)^2$$

$$\hat{y}_i = f_{\underline{w}}(x_i)$$

Polynomials of order *M* (red) fit to data constructed as t = sin(2πx) + ε (green)
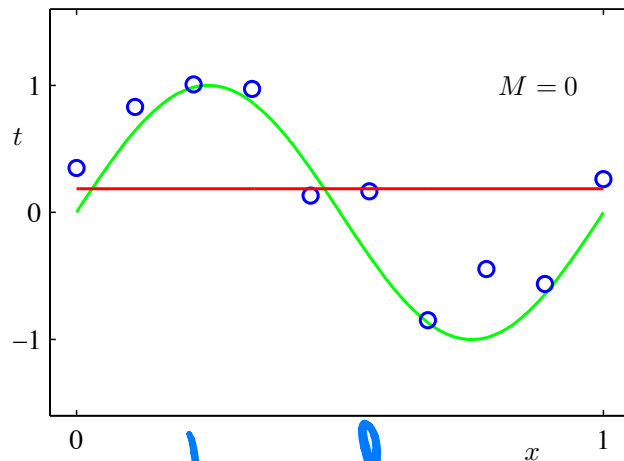
# P: Performance measure

Clustering



$$\text{within cluster sum of squares} = \sum_{i=1}^{n_\text{samples}} \min_{\mu_j \in C} \| \mu_j - x_i \|^2$$
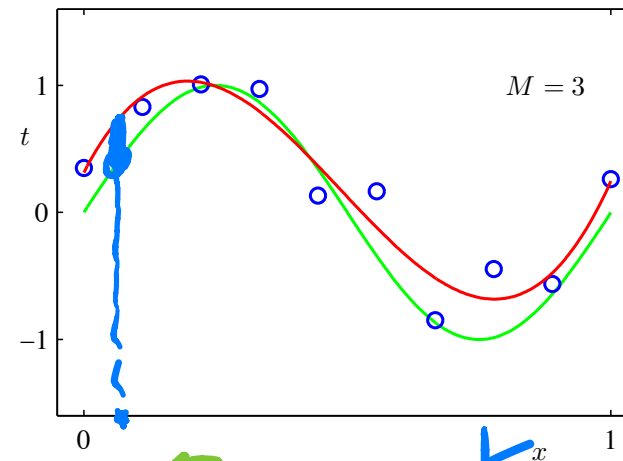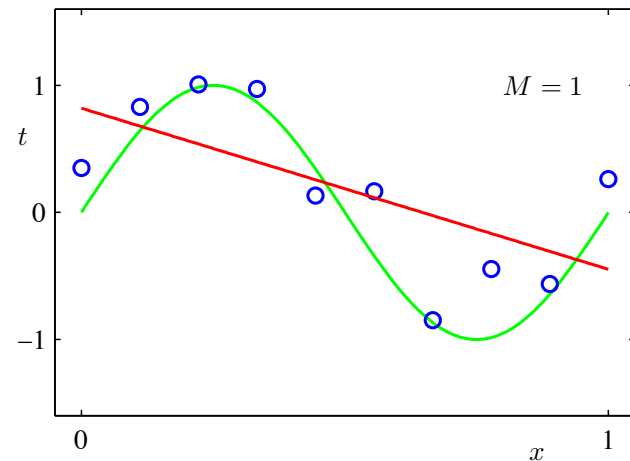
Expression matrix of genes (rows) for 64 human tumor samples (columns).  [source: ESL 1.3]
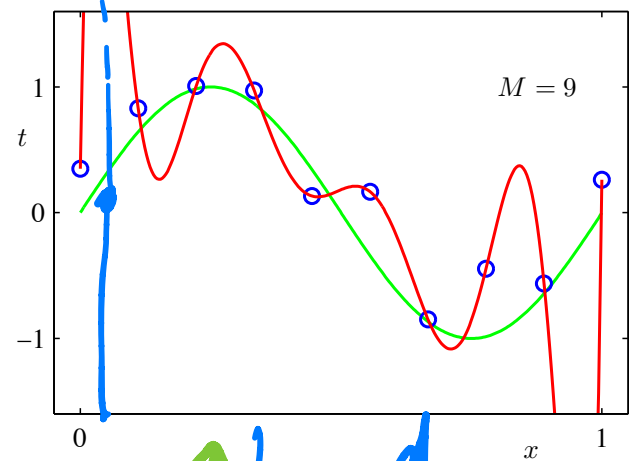
# P: Performance measure

$$\mathrm{MSE}(y, \hat{y}) = \frac{1}{n_\mathrm{samples}} \sum_{i=1}^{n_\mathrm{samples}} (y_i - \hat{y}_i)^2$$

*overfitting*



M = 0    M = 1    M = 3    M = 9

**bad**    **great**    **bad**

**Best performance on training set :**

**Best performance on new datapoints :**

Q: On which datapoints should performance be measured?

**Generalisation:**

*performance should be measured on new data (test data)*

# What is machine learning?

"A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T, as measured by P, improves with experience E."

- Tom M. Mitchell